

4 / 02 / 2020

## MACHINE LEARNING UDEMY

### LINEAR REGRESSION SCIKIT LEARN

- SUPERVISED LEARNING
- REGRESSION TOWARDS A MEAN
- DATA → LINEAR REGRESSION → MODEL

SUM OF SQUARED ERRORS

SUM OF ABSOLUTE ERRORS

SUM OF SQUARES OF THE RESIDUALS

X = FEATURES (NUMBER OF ROOMS)

Y = WHAT YOU WANT TO PREDICT (PRICE)

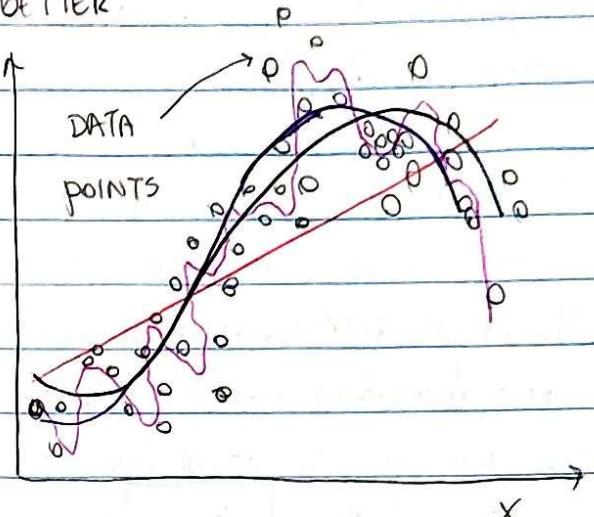
X-train, X-test, y-train, y-test

= train\_test\_split(x, y, test\_size  
= 0.4, random\_state=101)

### BIAS VARIANCE TRADE OFF

INCREASING THE COMPLEXITY OF A MODEL  
DOES NOT NECESSARILY MEAN THE MODEL

IS BETTER.



■ QUADRATIC

■ REAL CURVE

■ SPLINE

TRAIN TEST SPLIT : SPLITS DATA IN TWO,

ONE PART TO TRAIN THE MODEL AND THE

OTHER TO TEST IT. TEST SIZE IS THE PER-

CENTAGE OF TOTAL DATA THAT WILL BECOME

THE TEST DATA. RANDOM STATE IS HOW

SCIKIT LEARN SPLITS DATA.

SEE HOW, IN THIS CASE, THE MOST COMPLEX

MODEL ENCAPTURES A LOT OF DATA POINTS,

BUT IS FARAWAY FROM THE REAL CURVE

lm = LinearRegression()

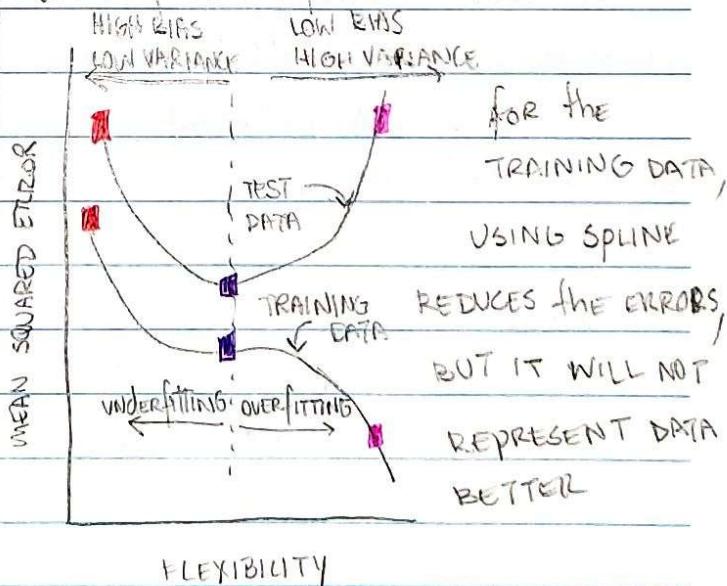
lm.fit(X-train, Y-train)

lm.coef\_

RECOMMENDED BOOK : INTRODUCTION TO  
STATISTICAL LEARNING BY GARRETH

JAMES, CHAPTERS 1 AND 2

(ISL BOOK)



# MACHINE LEARNING udemy

## LOGISTIC REGRESSION

INTRODUCTION TO STATISTICAL  
LEARNING, SECTIONS 4-4.3

• ERROR RATE (MISCLASSIFICATION RATE)

how often is it wrong?  $\frac{FP+FN}{TOTAL} = 9\%$

- CLASSIFICATION; DISCRETE CATEGORIES

SCIKIT LEARN

- SIGMOID FUNCTION

`sklearn.metrics.confusion_matrix`

$$\phi(z) = \frac{1}{1 + e^{-z}}$$

$$P = \frac{1}{1 + \exp(-(a + bx))}$$

→ VALUES RANGE FROM 0 TO 1

PRECISION: ABILITY of the classifier NOT TO

- CONFUSION MATRIX → MODEL EVALUAT.

LABEL AS positive A SAMPLE THAT

→ DESCRIBES the performance of  
the classification model

IS NEGATIVE:

$$\frac{TP}{TP+FP}$$

	PREDICTED FALSE	PREDICTED TRUE
ACTUAL FALSE	TN = 50	FP = 10
ACTUAL TRUE	FN = 5	TP = 100

RECALL: ABILITY of the classifier to find

ALL positives:  $\frac{TP}{TP+FN}$

TP = TRUE POSITIVE

PRECISION AND RECALL

TN = TRUE NEGATIVE

0 → WORST CASE, 1 → BEST

FP = FALSE POSITIVE = TYPE I ERROR

FN = FALSE NEGATIVE = TYPE II ERROR

SUPPORT: NUMBER of OCCURRENCES OF EACH  
CLASS IN y\_true (CORRECT VALUES)

- ACCURACY: how often is correct?

$$\frac{TP+TN}{TOTAL} = ACCURACY = 91\%$$

TOTAL

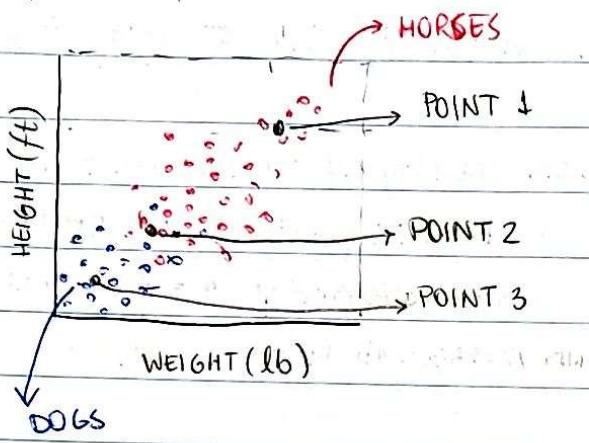
4 / 04 / 20

## MACHINE LEARNING (UDMII)

### K NEAREST NEIGHBORS (KNN)

- INTRODUCTION TO STATISTICAL LEARNING CHA
- CLASSIFICATION ALGORITHM BASED ON THE NEIGHBORING POINTS CALCULATING THE DISTANCE FROM  $x$  TO ALL POINTS

EXAMPLE: NEW POINTS: DOG OR HORSE?



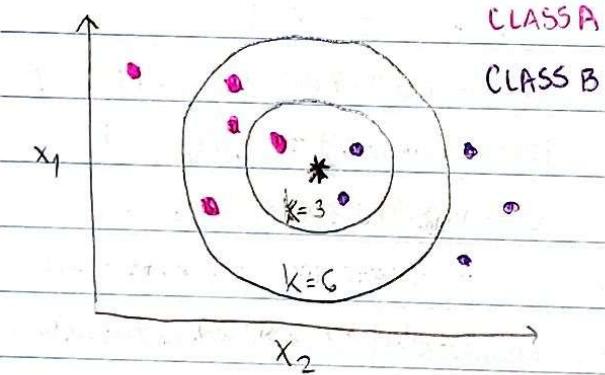
the greater the k number, more

BIAS, BETTER CUT off but ALSO MORE MISLABELED POINTS

SIMPLE MODEL WITH FEW PARAMETERS, BUT IT IS NOT RECOMMENDED FOR HIGH DIMENSIONAL DATA AND LARGE DATA SETS AND CATEGORICAL FEATURES

- SCALE DATA WITH `sklearn.preprocessing`.  
`STANDARDSCALER` TO `STANDARDIZE`
- FIND OPTIMAL K NUMBER by comparing ERROR RATES ( $y_{\text{test}} - \text{predictions}$ )

K NUMBER AFFECTS the CLASS A NEW POINT GOES INTO:



$k=3 \therefore$  NEW POINT IS FROM CLASS B

$k=6 \therefore$  NEW POINT IS FROM CLASS A

COMPARES HOW MANY NEIGHBORS THERE ARE

## MACHINE LEARNING UDEMY

### DECISION TREES AND RANDOM FORESTS

RANDOM FORESTS: COLLECTION OF DECISION TREES

→ MINIMIZATION OF ERRORS BY BIAS AND/OR VARIANCE

- INTRODUCTION TO STATS LEARN. CH8 → LIMIT OVERTFITTING

DECISION TREES: CLASSIFICATION & REGRESSION, NUMERICAL & CATEGORICAL

DIFFERENCES:

DECISION TREE

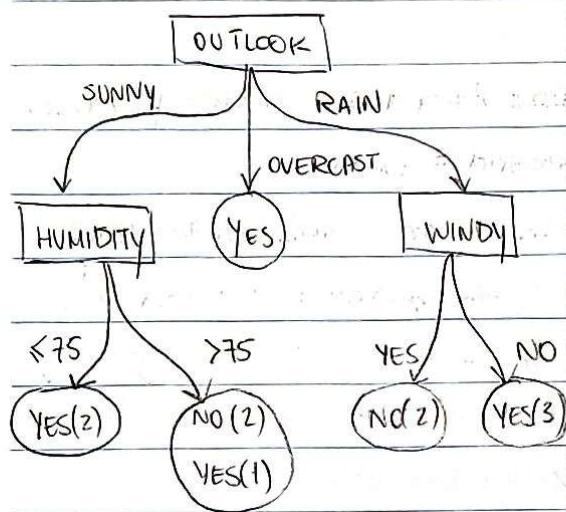
→ BUILT ON ENTIRE DATASET

USING ALL FEATURES

RANDOM FOREST

→ RANDOMLY SELECTS FEATURES

AND OBSERVATIONS



the RANDOM FOREST builds multiple decision trees with a number of features and observations and averages the results

DECISION TREE: `sklearn.tree.DecisionTreeClassifier`

RANDOM FOREST: `sklearn.ensemble.RandomForestClassifier`

`rfc = RandomForestClassifier (parameters)`

→ NODES (OUTLOOK, HUMIDITY, WINDY)

↳ n. estimators, max-features

→ EDGES (NODE OUTCOME)

number of trees to be used in the forest

max number of features to be considered

→ ROOT (OUTLOOK)

→ REGRESSION: USE n-FEATURES

→ LEAVES ("RESULT"; YES/NO)

→ CLASSIFICATION: USE  $\sqrt{n}$ -FEATURES

→ ENTROPY & INFORMATION GAIN

### DISADVANTAGES

- MAKES MOST OPTIMAL DECISION ON EACH STEP, but NOT THE GLOBAL OPTIMUM
- PRONE TO OVERTFITTING

4 / 12 / 20

## MACHINE LEARNING UDEMY

### SUPPORT VECTOR MACHINES (SVM)

```
from sklearn.model_selection import GridSearchCV
```

- ISL CHAPTER 9
- SUPERVISED LEARNING TO RECOGNIZE PATTERNS FOR CLASSIFICATIONS AND REGRESSION
- NON PROBABILISTIC BINARY LINEAR CLASSIFIER → CLASSIFICATION

param-grid =  
{"C": [0.1, 1, 10, 100, 1000],  
"gamma": [1, 0.1, 0.01, 0.001, 0.0001]}

grid = GridSearchCV (estimator, param-grid, verbose, refit)

→ estimator = SVC()

- MANY HYPERPLANES CAN EXIST BETWEEN CLASSES
- GOAL: CHOOSE A HYPERPLANE THAT MAXIMIZES THE MARGIN BETWEEN CLASSES

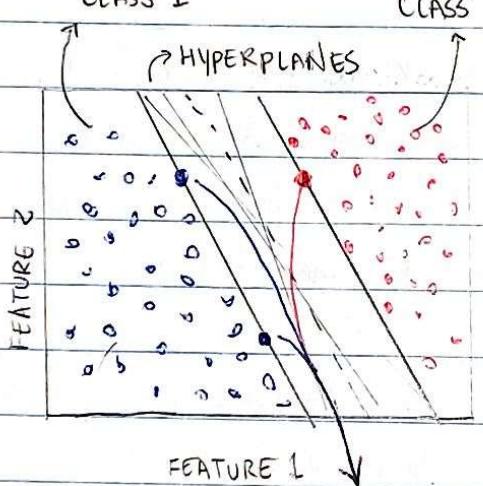
needs to provide a score function

→ param-grid = param-grid

DICTIONARY WITH PARAMETER NAMES

(i.e. C and gamma) AND VALUES

- (LIST OF SETTINGS TO BE TESTED)



C: CONTROLS THE COST OF MISCLASSIFICATION ON THE TRAINING DATA

↑C :: LOW BIAS, HIGH VARIANCE;  
HIGH PENALIZATION

GAMMA: FREE PARAMETER OF THE GAUSSIAN

RADIAL BASIS FUNCTION (kernel='rbf')

↑GAMMA :: HIGH BIAS, LOW VARIANCE;

SUPPORT VECTOR WON'T HAVE  
A WIDESPREAD INFLUENCE

→ verbose = 3 (default zero)

the higher the verbose, the more messages

→ refit = True

REFIT AN ESTIMATOR USING THE BEST FOUND  
PARAMETERS ON THE WHOLE DATASET

- THE VECTOR POINTS THAT TOUCH THE MARGIN LINES ARE THE SUPPORT VECTORS

- KERNEL TRICK → NON LINEARLY SEPARABLE DATA

- GRID SEARCH: GRID OF PARAMETERS USED TO TRY COMBINATIONS TO FIND GOOD SVM PARAMETERS

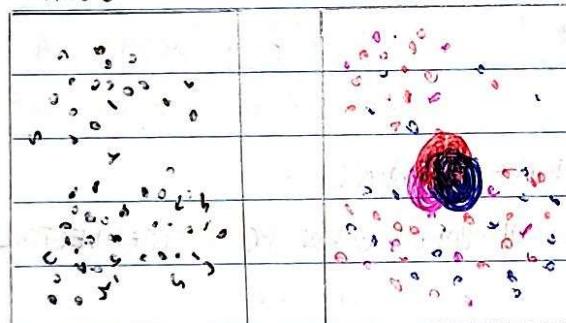
## MACHINE LEARNING UDEMY

### K MEANS CLUSTERING

- ISL CHAPTER 10
  - UNSUPERVISED LEARNING
  - DIVIDE DATA INTO DISTINCT GROUPS
  - ALGORITHM
- CHOOSE A NUMBER OF CLUSTERS  $K$
- RANDOMLY ASSIGN EACH POINT TO A CLUSTER
- ITERATE OVER THE COMPUTED CENTROIDS UNTIL SSE (SUM OF SQUARED ERRORS) DOES NOT PLUNGE WHEN COMPARED TO THE  $K$  NUMBER!

EXAMPLE: TAKE UNLABELED DATA AND ASSIGN EACH POINT TO A RANDOM CLUSTER. COMPUTE THE CENTROID.

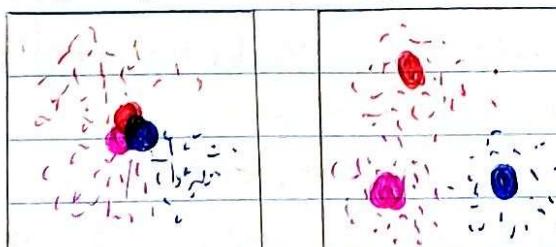
#### UNLABELED DATA



#### STEP 1

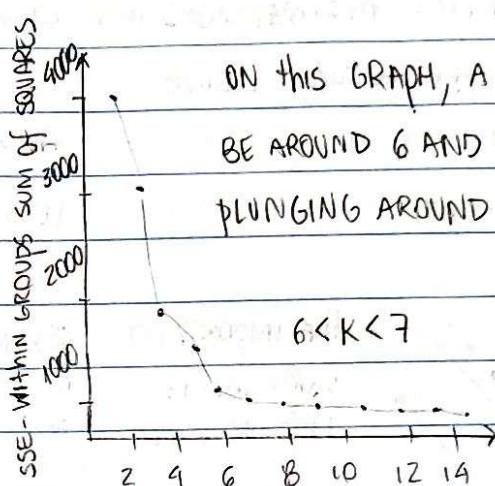
ASSIGN THE DATA POINTS TO THE CLUSTER FOR WHICH THE CENTROID IS CLOSEST

#### ITERATIONS



### CHOOSING K VALUE

- SSE:  $\sum (\text{POINT} - \text{CENTROID})^2$
- ↑ NUMBER of CLUSTERS  $\therefore \downarrow \text{SIZE of CLUSTERS}$
- CHOOSE A  $K$  AT WHICH THE SSE DECREASES ABRUPTLY (ELBOW EFFECT)



ON THIS GRAPH, A GOOD  $K$  WOULD BE AROUND 6 AND 7, AS IT STOPS PLUNGING AROUND THESE VALUES

$K = \text{NUMBER of CLUSTERS}$

4 13 20

## MACHINE LEARNING UDEMY

### PRINCIPAL COMPONENT ANALYSIS

- ISL CH 10.2
- UNSUPERVISED LEARNING
- COMPRESSION OF MANY DATA DIMENSIONS INTO A FEW THAT ARE THE MOST RELEVANT

MAIN IDEA: FLATTEN n-DIMENSIONS INTO A FEW WITHOUT LOSING VALUE OF DATA

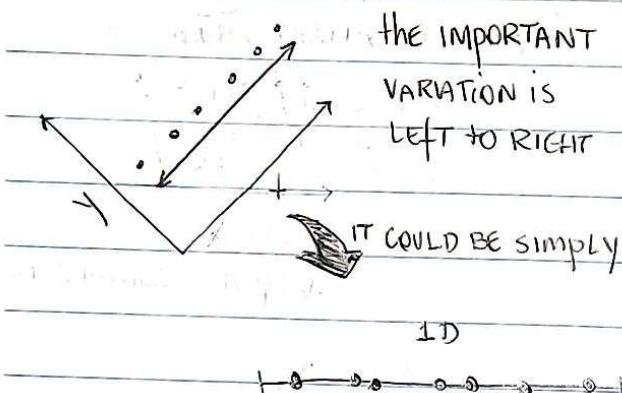
EXAMPLE: 2D DATA:

EXAMPLE: sklearn BREAST CANCER DATASET

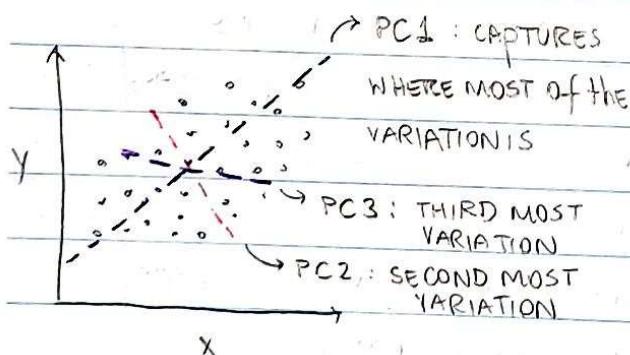
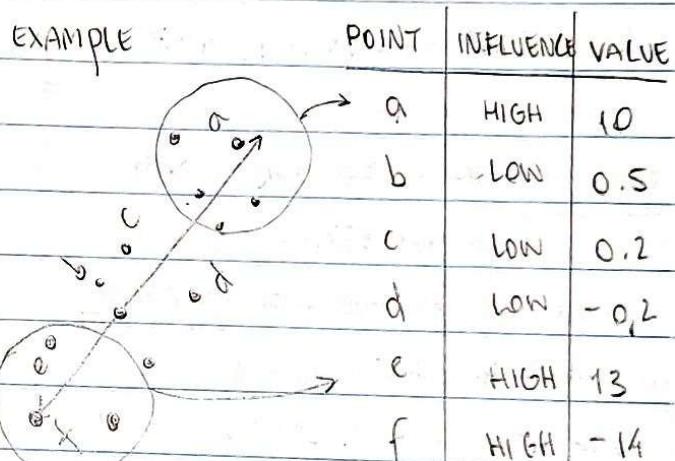
→ 30 FEATURES (RADIUS, AREA, SYMMETRY...)

CANCER DATA CAN BE TRANSFORMED WITH PCA TO HAVE ONLY TWO PRINCIPAL COMPONENTS. PC1 AND PC2 ARE COMBINATIONS OF WEIGHTS OF EACH OF THE 30 FEATURES

the closer the PC value is to zero, the less influence (weight) that point has



EXAMPLE: TAKE THE GREATEST VARIATIONS ON THE LINES THAT REPRESENT DATA



- THE PRINCIPAL COMPONENTS ARE A COMBINATION OF ALL OR SOME OF THE FEATURES

STEPS:

→ STANDARDIZE DATA

sklearn.preprocessing. StandardScaler  
scaler.fit(df) and scaler.transform(df)

→ PCA

sklearn.decomposition.PCA

pca.fit(scaledData) and pca.transform

## MACHINE LEARNING UDEMY

### RECOMMENDER SYSTEMS

→ RECOMMENDER SYSTEMS by JANNACH AND ZANKER

→ COLLABORATIVE FILTERING (CF)

RECOMMENDATIONS BASED ON USERS

RATING; "AUTO" FEATURE LEARNING

→ CONTENT-BASED

RECOMMENDATIONS BASED ON

ATTRIBUTES OR SIMILARITIES OF ITEMS

• TF-IDF: TERM FREQUENCY - INVERSE DOCUMENT FREQUENCY

→ COUNT HOW MANY TIMES A WORD OCCUR

→ TERM FREQUENCY (TF): IMPORTANCE OF TERM t IN DOCUMENT d ∴  $TF = f(d, t)$

→ WEIGHTS COUNTS SO FREQUENT TOKENS GET LOW WEIGHT

→ INVERSE DOCUMENT FREQUENCY (IDF): IMPORTANCE OF TERM t IN CORPUS (GROUP OF ALL DOCS)

∴  $IDF(t) = \log(D/t)$

D: NUMBER OF DOCUMENTS WITH TERM

D: total NUMBER OF DOCUMENTS

### NATURAL LANGUAGE PROCESSING

• COMPILE AND FEATURIZE DOCUMENTS → TF-IDF (W): TERM X WITHIN DOCUMENT Y

EXAMPLE: DOCS "BLUE HOUSE" AND  
"RED HOUSE"

→ FEATURIZE BY WORD COUNT  
(red, blue, house)

"BLUE HOUSE" → (0, 1, 1)

"RED HOUSE" → (1, 0, 1)

$$W_{x,y} = t_{x,y} \cdot \log\left(\frac{N}{df_x}\right) \rightarrow \text{total of docs}$$

↓  
NUMBER OF DOCS CONTAINING X  
↓  
FREQUENCY OF X IN Y

this CALCULATION CONSIDERS HOW IMPORTANT  
A WORD IS BESESIDES ITS COUNT

• BAG OF WORDS: A DOCUMENT REPRESENTED  
AS A VECTOR

WORD/DOCUMENT SIMILARITY

$$\text{SIMILARITY}(A, B) = \cos(\theta) = \frac{A \cdot B}{\|A\| \cdot \|B\|}$$

• BAG OF WORDS CAN BE IMPROVED BY ADJUSTING  
WORD COUNTS BASED ON THEIR FREQUENCY  
IN CORPUS (GROUP OF ALL DOCUMENTS)

• STOP WORDS: FREQUENT WORDS USED IN A  
LANGUAGE; UNIMPORTANT TO NLP

• TOKENIZATION: CONVERSION OF PURE TEXT INTO  
TOKENS (THE IMPORTANT PART OF THE TEXT)  
BY REMOVING STOP WORDS AND PUNCTUATION, i.e.

• STEMMING AND LEMMATIZATION

the boy's cars are different colors

the boy car be differ color

4 / 24 / 2020

## MACHINE LEARNING UDACITY

### NEURAL NETS AND DEEP LEARNING

WEIGHT ( $w$ ): ADJUSTS INPUT  $x_i$  STRENGTH

TOPICS:

- PERCEPTRON MODEL
- ACTIVATION FUNCTIONS
- COST FUNCTIONS
- FEED FORWARD NETWORKS
- BACK PROPAGATION

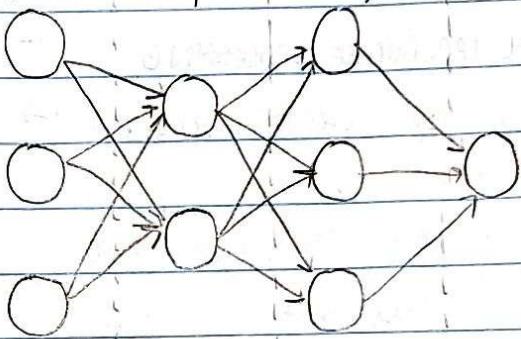
CODING TOPICS:

- TENSORFLOW 2.0 (KERAS)
- ARTIFICIAL NEURAL NETWORKS
  - ↳ REGRESSION
  - ↳ CLASSIFICATION

BIAS ( $b$ ): OFFSET TO MAKE  $x_i * w_i$  REACH A CERTAIN THRESHOLD BEFORE HAVING AN EFFECT; SHIFTS THE ACTIVATION FUNCTION.

### MULTILAYER PERCEPTRON MODEL

INPUT	HIDDEN LAYERS	OUTPUT
LAYER 1	LAYER 2	LAYER 3
		LAYER 4



### PERCEPTRON MODEL

• OUTPUT OF A LAYER IS THE INPUT OF

ANOTHER

IDEA:

INPUT      FUNCTION      OUTPUT

$$x_1 * w_1 + b$$

$y$

↳ INTERACTION AND RELATIONSHIPS

BETWEEN FEATURES

• DEEP NEURAL NETWORK: IF THERE ARE 2 OR MORE HIDDEN LAYERS

• UNIVERSAL APPROXIMATION THEOREM:

ANY CONVEX CONTINUOUS FUNCTION CAN BE APPROXIMATED BY A FEED-FORWARD NETWORK WITH A SINGLE HIDDEN LAYER CONTAINING A FINITE NUMBER OF NEURONS

$$\hat{y} = \sum_{i=1}^n x_i w_i + b_i$$

WHERE:  $x_i$  = INPUT,  $w_i$  = WEIGHT

$$b_i = \text{BIAS}; \hat{y} = \text{OUTPUT}$$

## MACHINE LEARNING UDEMY - NEURAL NETS AND DEEP LEARNING

### ACTIVATION FUNCTIONS

- they set boundaries to output values from the neuron

$$z = x * w + b$$

$x \rightarrow$  INPUT

$w \rightarrow$  WEIGHT (INPUT STRENGTH)

$b \rightarrow$  BIAS (offset)

i.e. if  $b = -10$ : the effects of  $x * w$  won't start

to overcome the bias until their product is greater than 10

ACTIVATION FUNCTIONS LIMITS  $z$

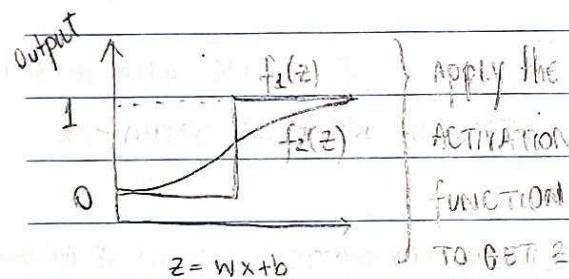
### TABLE OF COMMON ACTIVATION FUNCTIONS

NAME	PLOT	EQUATION
BINARY STEP		$\sigma(z) = \begin{cases} 0, z < 0 \\ 1, z \geq 0 \end{cases}$ {0,1}
SIGMOID		$\sigma(z) = \frac{1}{1 + \exp(-z)}$ (0,1)
tanh		$\sigma(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)}$ (-1,1)
RECTIFIED LINEAR UNIT		$\sigma(z) = \begin{cases} 0, z \leq 0 \\ z, z > 0 \end{cases}$ [0, infinity)
SOFTMAX		PROBABILITY DISTRIBUTION OF EVENTS OVER K EVENTS. PROBABILITY OF TARGET CLASS OVER ALL CLASSES

EXAMPLE: CLASSIFICATION problem

(OUTPUT IS EITHER 0 OR 1)

MUTUALLY EXCLUSIVE CLASS: i.e. SOFTMAX



	RED	GREEN	BLUE
DATA POINT 1	1	0	0
DATA POINT 2	0	1	0
...	...	...	...
DATA POINT N	0	0	1

### MULTICLASS CLASSIFICATION

#### NON EXCLUSIVE CLASS

↳ A data point can have multiple classes/categories (i.e. photo tags)

#### MUTUALLY EXCLUSIVE CLASSES

↳ ONLY ONE CLASS PER DATA

POINT (i.e. photo is B&W or color)

NON-EXCLUSIVE CLASSES: i.e. SIGMOID

	A	B	C	
DATA 1	1	1	0	CLASS 1: 0.8
DATA 2	1	0	0	CLASS 2: 0.3
DATA 3	0	1	1	CLASS 3: 0.5
...	...	...	...	...
DATA N	0	1	0	

4 / 24 / 2020

## MACHINE LEARNING UDEMY - NEURAL NETS AND DEEP LEARNING

### LOSS FUNCTIONS & GRADIENT DESCENT

$$\text{COST} = C(W, B, S^r, E^r)$$

WHERE  $C = \text{COST}$ ,  $W = \text{WEIGHT}$ ,  $B = \text{BIAS}$

$S^r = \text{INPUT OF SINGLE TRAINING}$

SAMPLE

$E^r = \text{DESIRED OUTPUT OF TRAINING}$

SAMPLE

FORMULAS: QUADRATIC COST FUNCTION

$$C = \frac{1}{2n} \sum_x \|y(x) - a^L(x)\|^2$$

$$a = \sigma(z)$$

$$z = Wx + b$$

WHERE:  $y(x)$ : TRUE VALUE

$a^L(x)$ : NEURON'S PREDICTION OR

ACTIVATION OUTPUT FOR LAYER

L (LAST LAYER)

$\sigma(z)$ : ACTIVATION FUNCTION

n: NUMBER OF POINTS

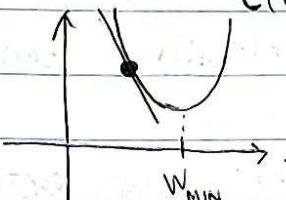
GOAL: MINIMIZE THE COST (IT'S THE DIFFERENCE

BETWEEN THE REAL OUTPUT AND THE PREDICTION)  $\therefore$  DERIVATIVE AND SOLVE FOR ZERO  $\therefore$  GRADIENT DESCENT

SIMPLIFIED EXAMPLE: CALCULATE SLOPE

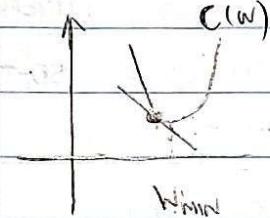
$C(W)$  AND MOVE IT

DOWNTWARD UNTIL CONVERGE TO ZERO



STEP SIZE = LEARNING RATE. IF IT ISN'T CONSTANT, THEN IT'LL BE AN ADAPTIVE

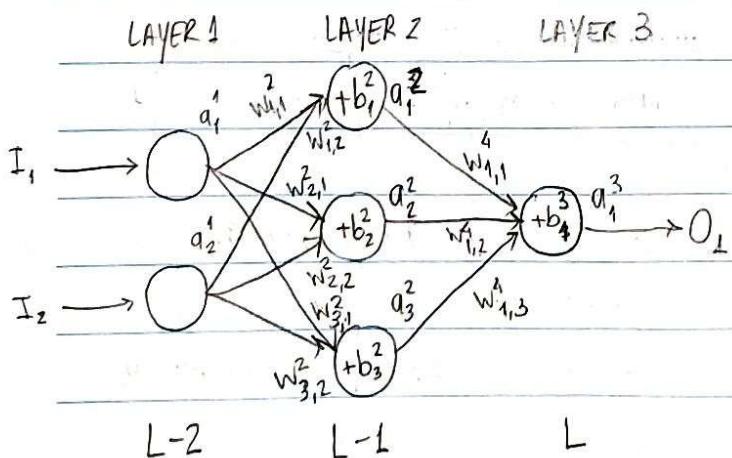
GRADIENT DESCENT ( $\nabla C$ )



ADAM: METHOD FOR STOCHASTIC OPTIMIZATION.

VARIATION OF RMSPROP: OPTIMIZER

STOCHASTIC APPROXIMATION: ITERATIVE METHODS  
FOR ROOT-FINDING OR OPTIMIZATION



DERIVATIVE  $\rightarrow$  VECTORS

GRADIENTS  $\rightarrow$  TENSORS

FOR CLASSIFICATION PROBLEMS: CROSS ENTROPY  
(MEASURE OF DIFFERENCE BETWEEN TWO  
PROBABILITY DISTRIBUTIONS FOR EVENTS)

## MACHINE LEARNING UDACITY - NEURAL NETS AND DEEP LEARNING

### BACKPROPAGATION

- EVALUATES THE EXPRESSION FOR THE DERIVATIVE OF THE COST FUNCTION AS A PRODUCT OF DERIVATIVES BETWEEN EACH LAYER BACKWARDS WITH THE GRADIENT OF THE WEIGHTS BETWEEN EACH LAYER BEING A SIMPLE MODIFICATION OF THE PARTIAL PRODUCTS (THE BACKWARDS PROPAGATED ERROR)

EXAMPLE:  $C(w_1, b_1, w_2, b_2)$

$$w_1 + b_1 \quad w_2 + b_2$$

L-2

L-1

$$z = wx + b$$

ACTIVATION FUNCTION  $a = \sigma(z)$

$$z^L = w^L a^{L-1} + b^L$$

$$a^L = \sigma(z^L)$$

$$C_0 = (a^L - y)^2$$

HOW SENSITIVE IS THE COST FUNCTION

CHANGE WITH RESPECT TO WEIGHT/BIAS?

$$\frac{\partial C_0}{\partial w^L} \text{ and } \frac{\partial C_0}{\partial b^L}$$

GOAL: MINIMIZE THE OUTPUT OF THE ERROR VECTOR ON LAYER L BY ADJUSTING WEIGHT AND BIAS GOING BACKTHROUGH THE NETWORK USING THE GRADIENT

$$\text{ERROR VECTOR: } \delta^L = \nabla_a C \odot \sigma'(z^L)$$

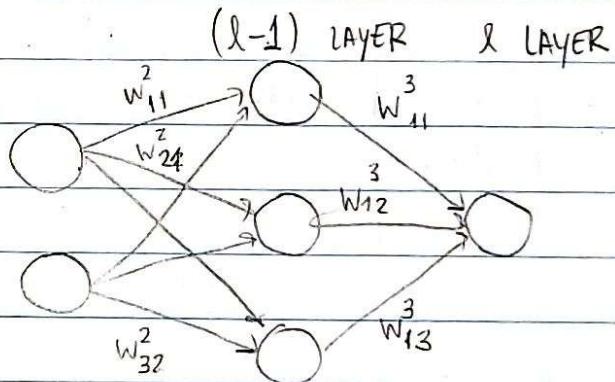
WHERE L = LAST LAYER

l = CURRENT MIDDLE LAYER

• = HADAMARD PRODUCT

$$\nabla_a C = (a^L - y)$$

BACK TO THE NEURAL NETWORK EXAMPLE:



NOTATION:  $w_{jk}^l, b_j^l, a_j^l$

j: NEURON IN THE  $l^{\text{th}}$  LAYER

k: NEURON IN THE  $(l-1)^{\text{th}}$  LAYER

l: LAYER

$w_{jk}^l$ : WEIGHT FROM k TO j LAYER

$b_j^l$ : BIAS OF NEURON j IN LAYER l

$a_j^l$ : ACTIVATION OF NEURON j IN LAYER l

4 / 25 / 2020

## MACHINE LEARNING UDEMY - NEURAL NETS AND DEEP LEARNING

$$a^l = \sigma \left( \sum_k w_{jk}^l \cdot a_k^{l-1} + b_j^l \right)$$

ACTIVATION ON LAYER  $l$  IS BASED ON THE ACTIVATION OF THE PREVIOUS ( $l-1$ ) LAYER

BACKPROPAGATION FUNDAMENTAL EQUATIONS

BP1: ERROR IN THE OUTPUT LAYER

$$\delta_j^l = \frac{\partial C}{\partial a_j^l} \sigma'(z_j^l)$$

MATRIX FORM:

$$a^l = \sigma (w^l a^{l-1} + b^l)$$

$$\text{IF } C = \frac{1}{2} \sum_j (y_j - a_j^l)^2$$

$$\text{THEN } \frac{\partial C}{\partial a_j^l} = (a_j^l - y_j)$$

$$\delta^l = \nabla_a C \odot \sigma'(z^l) \text{ AS MATRIX form}$$

WHERE:

$$w^l = \begin{bmatrix} k_1 & k_2 & \dots & k_n \\ j_1 & \vdots & \ddots & \vdots \\ j_2 & & \ddots & \vdots \\ j_n & & & \ddots \end{bmatrix}$$

WITH SIMPLIFIED COST:

$$\delta^l = (a^l - y) \odot \sigma'(z^l)$$

BP2: ERROR FROM ERROR IN THE NEXT LAYER ( $\delta^{l+1}$ )

$$a^l = \begin{bmatrix} k_1 & \vdots & w_{jk}^l \cdot a_{k1}^l \\ k_2 & \vdots & \\ k_n & \vdots & \end{bmatrix}$$

$$\delta^l = ((W^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l)$$

TRANSPOSED: MOVES THE ERROR BACKWARDS  
THROUGH THE NETWORK

$$b^l = \begin{bmatrix} j_1 \\ j_2 \\ \vdots \\ j_n \end{bmatrix}$$

BP3: CHANGE OF COST WITH BIAS

$$z_j^l = \sum_k w_{jk}^l a_k^{l-1} + b_j^l$$

$$z^l = W a^{l-1} + b^l$$

$$\frac{\partial C}{\partial b_j^l} = \delta_j^l$$

$z$  = WEIGHTED INPUT

BP4: CHANGE OF COST WITH WEIGHT

HADAMARD PRODUCT:  $s \odot t$

$$\begin{bmatrix} 1 \\ 2 \end{bmatrix} \odot \begin{bmatrix} 3 \\ 4 \end{bmatrix} = \begin{bmatrix} 1 \cdot 3 \\ 2 \cdot 4 \end{bmatrix} = \begin{bmatrix} 3 \\ 8 \end{bmatrix}$$

$$\frac{\partial C}{\partial w_{jk}^l} = a_{k1}^{l-1} \delta_j^l$$



## BACKPROPAGATION ALGORITHM:

① INPUT  $x$ : SET THE CORRESPONDING ACTIVATION  $a$  FOR THE INPUT

LAYER.  $a$  FEEDS NEXT LAYER AS  $a^{l-1}$ .

② FEEDFORWARD: FOR EACH LAYER

$l = 2, 3, \dots, L$ , COMPUTE

$$z^l = w^l a^{l-1} + b^l \text{ and } a^l = \sigma(z^l)$$

③ OUTPUT ERROR  $\delta^L$ : COMPUTE VECTOR

$$\delta^L = \nabla_a C \odot \sigma'(z^L)$$

$a - y$  (BP)

④ BACKPROPAGATE THE ERROR: FOR EACH

LAYER  $l = L-1, L-2, \dots, 2$ , COMPUTE

$$\delta^l = ((w^{l+1})^\top \delta^{l+1}) \odot \sigma'(z^l)$$

⑤ OUTPUT: GRADIENT COSTS

$$\frac{\partial C}{\partial w_{jk}} = a_k^{l-1} \delta_j^l \text{ and } \frac{\partial C}{\partial b_j} = \delta_j^l$$

$a_k^{l-1}$

$\delta_j^l$

4 / 27 / 2020

## MACHINE LEARNING UDEMY

### BIG DATA AND PYSPARK

from pyspark import SparkContext

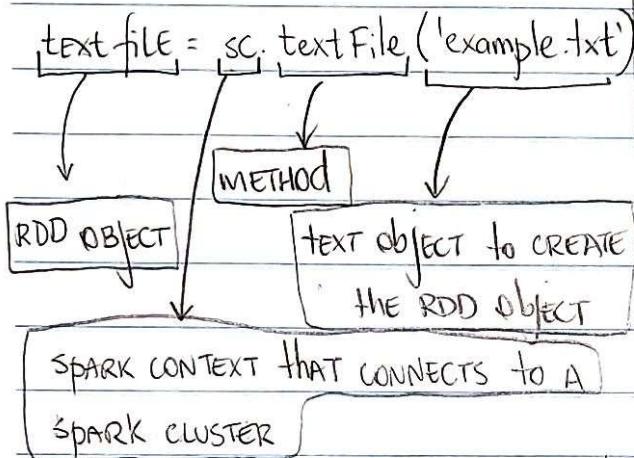
sc = SparkContext()

%% writefile example.txt

first line

second line

third line



RDD: ACTIONS: RETURNS VALUES

TRANSFORMATIONS: RETURNS POINTERS  
to NEW RDD's

textfile.count() → 3 (lines)

textfile.first() → 'first line'

secfind = textfile.filter(lambda

line: 'SECOND' in line)

secfind.collect() → ['second line']

secfind.count() → 1

RDD: RESILIENT DISTRIBUTED DATASET

TRANSFORMATION: SPARK OPERATION THAT PRODUCES AN RDD

ACTION: SPARK OPERATION THAT PRODUCES A LOCAL OBJECT

SPARK JOB: SEQUENCE OF TRANSFORMATIONS WITH A FINAL ACTION

COMMON RDD TRANSFORMATIONS:

- FILTER()
- UNION()
- MAP()
- DISTINCT()
- FLATMAP()
- SORTBY()
- SAMPLE()
- REDUCEBYKEY

COMMON RDD ACTIONS:

- COLLECT()
- SUM()
- TAKE()
- MEAN()
- TOP()
- STDEV()
- TAKESAMPLE()